



# **Sentence Encoder Assembly for Ad-hoc Video Search**

**Fangming Zhou, Aozhu Chen, Xirong Li**

**AI & Media Computing Lab, Renmin University of China**

TRECVID 2020 Workshop

2020-12-8



# Ad-hoc Video Search

How to retrieve unlabeled videos for ad-hoc textual queries?



## Using **cross-modal representation learning**

- Sentence representation
- Video representation
- Common space



# What We Focus This Year

## How to fully utilize multiple text encoders?

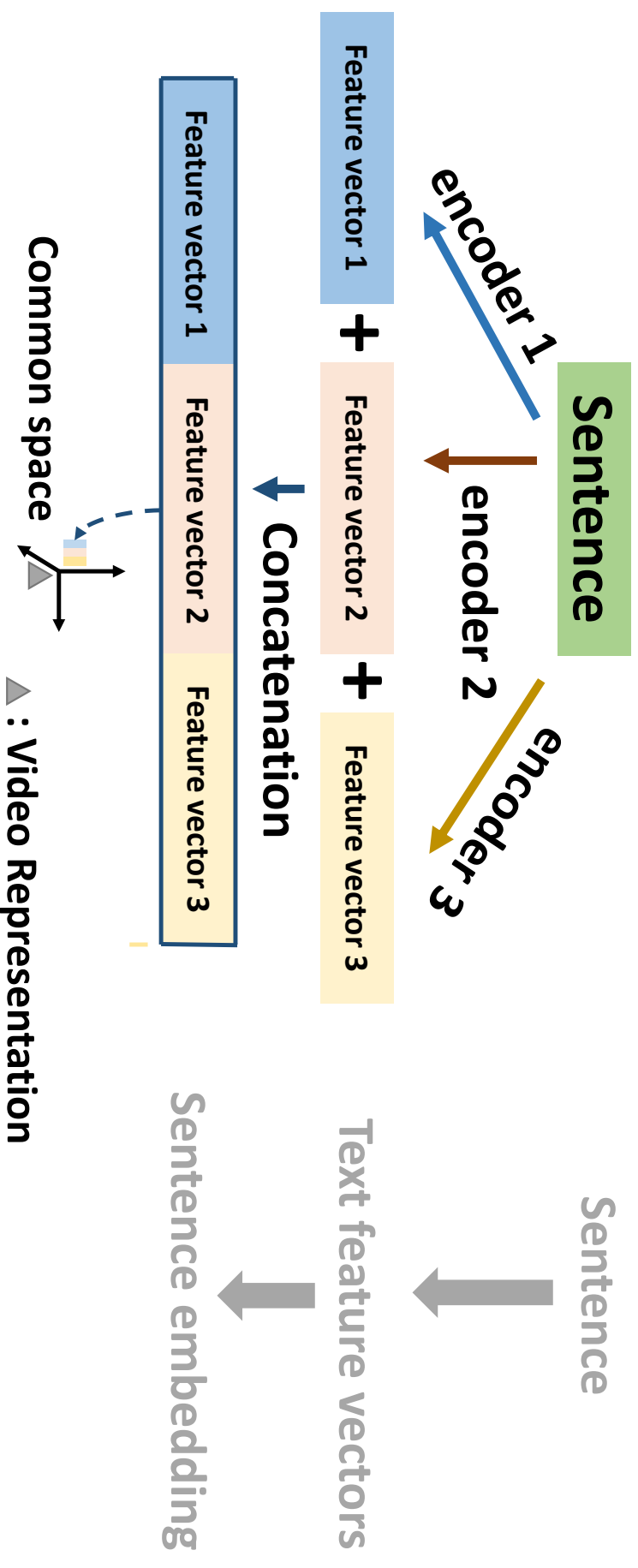
Previously some works combine multiple sentence encoders together as the sentence embedding

Works	Sentence encoders used
W2VV++ [Li et al., ACMMM'19] (AVS 2018 winner)	Bag-of-word, word2vec, GRU
Dual Encoding [Dong et al., CVPR'19]	Bag-of-word, bi-GRU, 1d-CNN
.....	.....

Li et al., W2VV++: Fully deep learning for ad-hoc video search, ACMMM 2019  
Dong et al., Dual Encoding for Zero-Example Video Retrieval, CVPR 2019

# What We Focus This Year

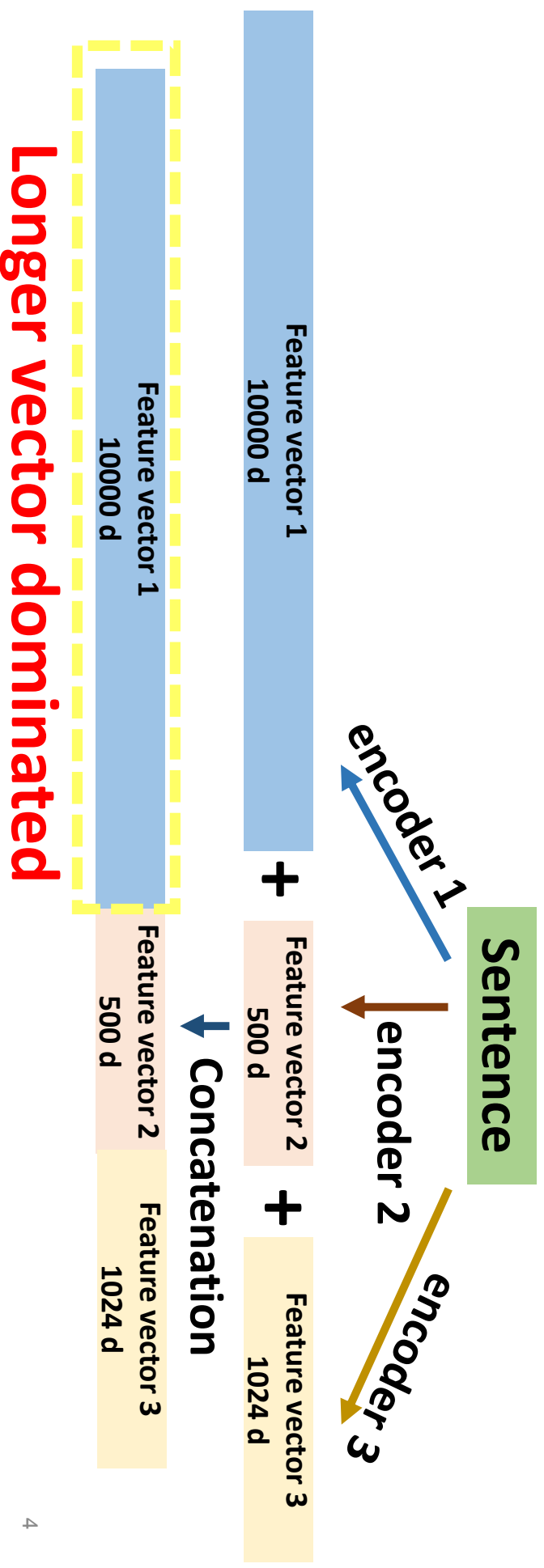
But they just concatenate text feature vectors



# What We Focus This Year

Their **disadvantages**? **Longer vector is dominant**

- The combined feature can be dominated by a specific encoder



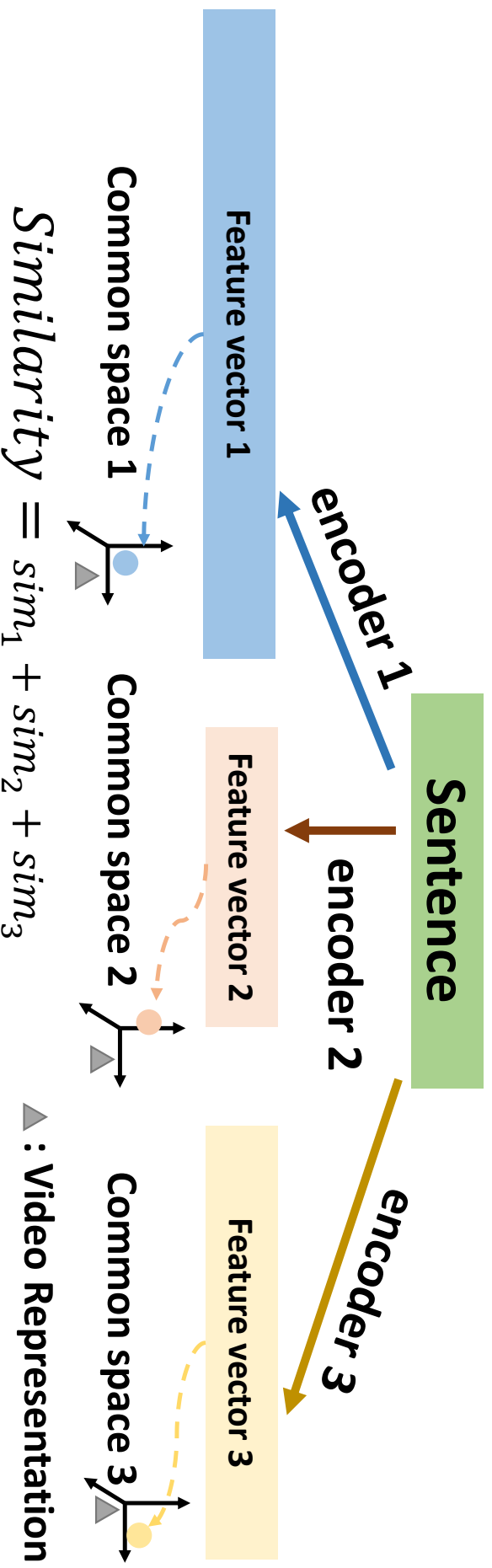
# How We Do

## Sentence Encoder Assembly (SEA)

A new and general architecture

### ① Multi-space Learning

Long vector dominated problem is solved



$$\text{Similarity} = \text{sim}_1 + \text{sim}_2 + \text{sim}_3$$

▲: Video Representation

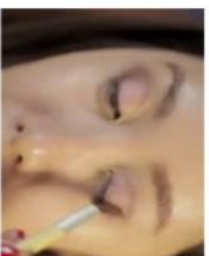
# How We Do

## Loss for multi space learning

We use the marginal ranking loss[1] to select the hard negative examples during training

- Selection based on the **combined similarity (a)**
- Selections based on **individual similarities (b-d)**

A female giving a nail art tutorial



(a)



(b)



(c)



(d)

More diverse hard negatives should be used

*Combined*

*BoW*

*w2v*

*GRU*

[1] Faghri et al., VSE++: Improving visual-semantic embeddings with hard negatives, BMVC 2018

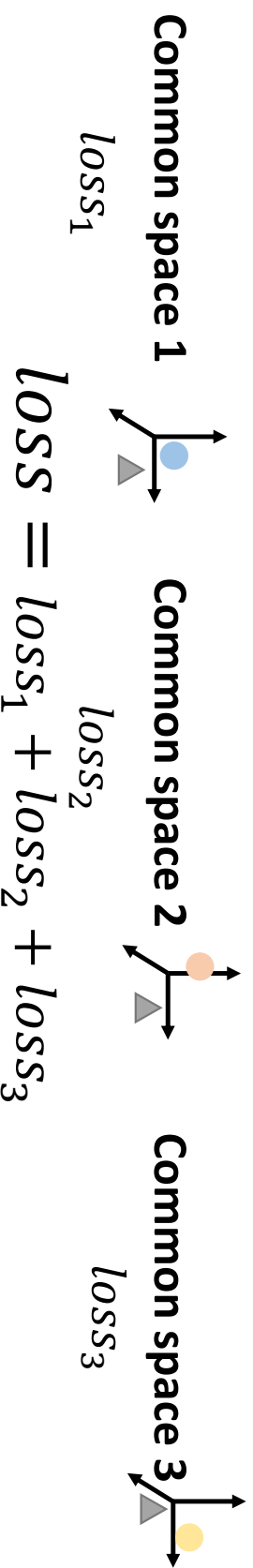


# How We Do

## Sentence Encoder Assembly (SEA)

### ② Multi-loss Learning

- Learning  $k$  common spaces for  $k$  encoders
- **Combined loss**  $loss = \sum_{i=1}^k loss_i(\text{sentence}, \text{video})$
- 30% extra hard negatives can be provided







# Is SEA Effective?

Retrospective experiments on TV16-19 in [1] can prove:

Sentence encoders	Model	TV16	TV17	TV18	TV19	SUM
bow,w2v	W2VV++	0.144	0.218	0.111	0.143	0.616
	SEA	0.157	0.234	0.128	0.166	0.685
bow,w2v,gru	W2VV++	0.162	0.223	0.101	0.139	0.625
	SEA	0.150	0.234	0.122	0.166	0.672
bow,w2v,bigru	W2VV++	0.161	0.217	0.104	0.135	0.617
	SEA	0.164	0.228	0.125	0.167	0.684
bow, w2v,bert	W2VV++	0.151	0.225	0.102	0.128	0.606
	SEA	0.153	0.228	0.121	0.148	0.650
bow,w2v,gru,bert	W2VV++	0.143	0.193	0.093	0.101	0.530
	SEA	0.160	0.231	0.121	0.154	0.666
bow,w2v,bigru,bert	W2VV++	0.158	0.206	0.090	0.105	0.559
	SEA	0.159	0.229	0.117	0.155	0.660

- Darker color indicates higher infAP
- SEA models surpass almost all the corresponding W2VV++ models

**Yes, it is**

\* W2VV++ uses a single common space

[1] Li et al., SEA: Sentence encoder assembly for video retrieval by textual queries, T-MM 2021.



# Choice of Sentence Encoder

We consider 5 sentence encoders

1. Bag-of-word (BoW)
2. word2vec ( $w2v$ )
3. NetVlad
4. bi-GRU \*
5. BERT \*

Different combinations are designed

1. SEA (BoW, NetVlad)
2. SEA (BoW,  $w2v$ )
3. SEA (BoW,  $w2v$ , bi-GRU)
4. SEA (BoW,  $w2v$ , bi-GRU, BERT)

↑  
complexity

\* Indicates **Sequence model**



# Datasets and Visual Features

Datasets	Usage
msrvtt10k	training
tgif	training
TV2016 VTT training set	validation
the Google's Conceptual Captions (GCC)	pre-training

**Image caption dataset pre-training is not suitable for models with C3D feature**

Video Features	Dim.
ResNext-101 (frame-level)	2,048
ResNet-152 (frame-level)	2,048
C3D (video-level)	2,048



Combined Video Features	Dim.
ResNext-101 + ResNet-152	4,096
ResNext-101 + ResNet-152 + C3D	6,144

**We use combined video features for better performance**



# Submissions

(fully automatic track)

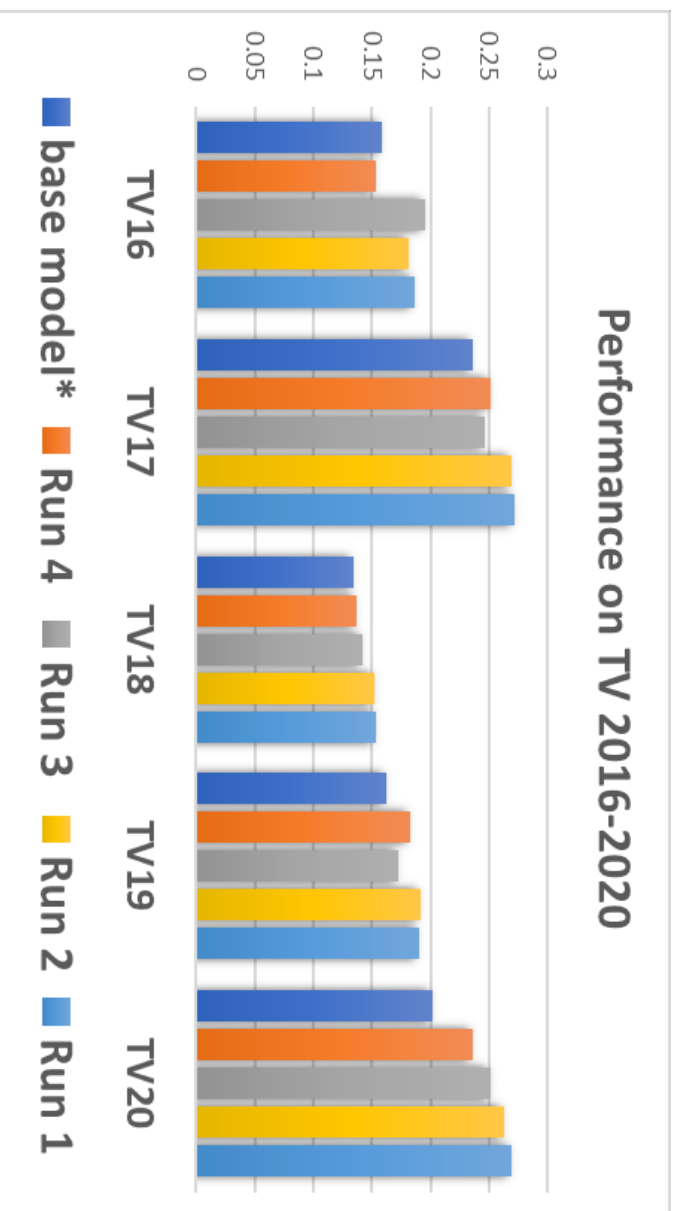
We submitted the following 4 runs:

SEA model with different setups and their combinations

Run id	Description
Run 4	SEA (BoW, w2v), <b>ResNeXt-ResNet-C3D</b>
Run 3	SEA (BoW, NetVlad), ResNeXt-ResNet, and <b>pre-trained on GCC</b>
Run 2	Late average fusion of three SEA models
Run 1 (primary run)	Late average fusion of four SEA models



# Performance on TV2016-2020 AVS Task



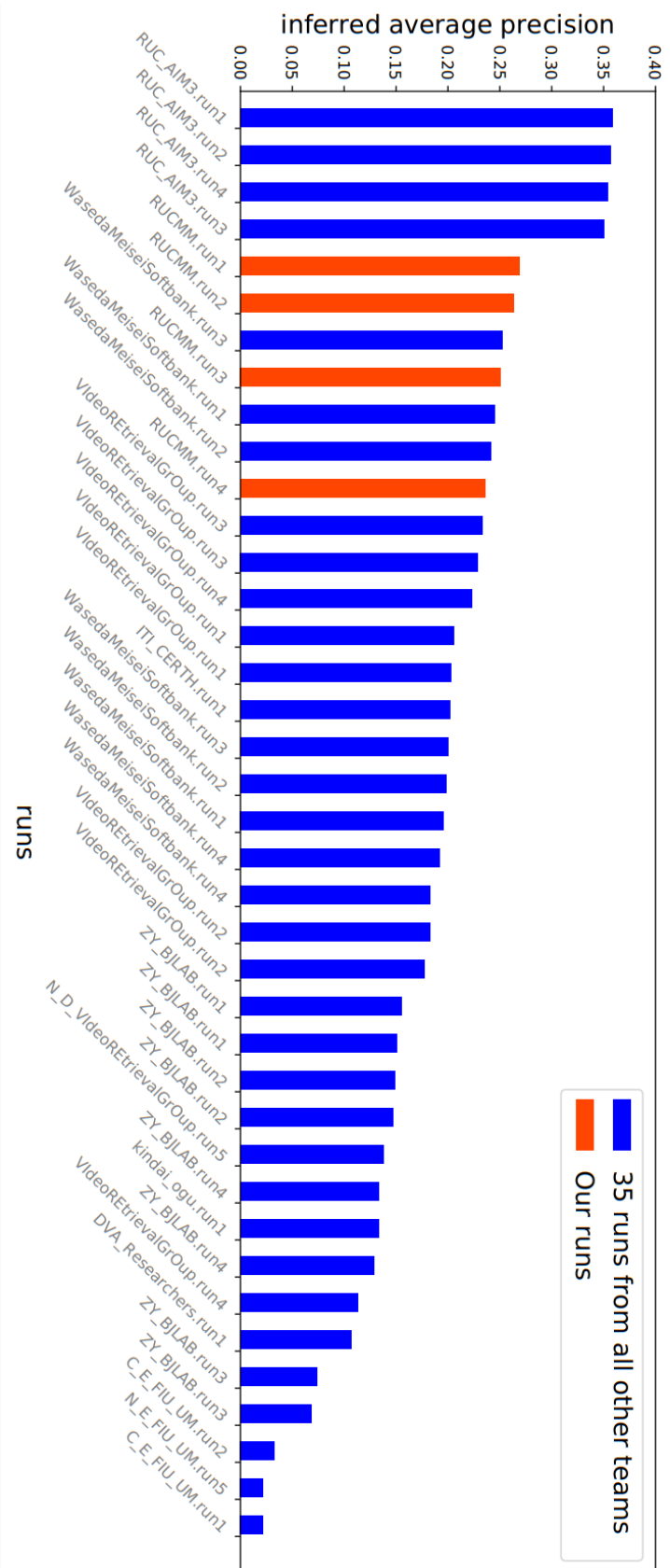
- Appending extra C3D feature helps (Run 4)
- Pre-training strategy helps (Run 3)
- Model ensemble helps (Run 2, Run 1)

\* **base model:** SEA(Bow, NetVlad), no C3D feature included, no pre-training



# All fully automatic AVS submissions

Our submissions ranked the 2nd





# Results of individual topics

More **blue** is better

More **red** is worse

Can be divided into 3 types:

- Easy topics: all **blue**
- Not easy topics: not all **blue**
- Hard topics: all **red**

topic id	runs1	runs2	runs3	runs4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455

# Case Study

Easy query

- All models perform well

656 a long haired man

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.312
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



shot04236\_8\_43



shot06870\_67\_23



shot06638\_91\_23



shot02033\_106\_212



shot06870\_69\_57



shot06638\_37\_11



shot06638\_24\_11



shot06638\_64\_11



shot01802\_55\_12



shot06638\_10\_103



# Case Study

## Easy query

- All models perform well

642 a person paddling kayak in the water

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
<b>642</b>	<b>0.522</b>	<b>0.541</b>	<b>0.454</b>	<b>0.516</b>
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.058
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



shot06958\_222\_150

shot02312\_93\_383

shot02312\_94\_29

shot06958\_159\_130

shot02862\_131\_45

# Case Study

No Easy query

- Not all models perform well

644 sailboats in the water

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.155
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



Run3 works well

# Case Study

No Easy query

- Not all models perform well

644 sailboats in the water

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.151
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



Run4 not works well, only shows water and boat



# Case Study

No Easy query

- Not all models perform well

644 **sailboats** in the water

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.000	0.000
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.199	0.194
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455

Datasets

Usage

msrvtt10k + tgif

training

7 captions have 'sailboat'

GCC

pre-training

853 captions have 'sailboat'

**Pre-training** makes the difference



# Case Study

## Hard query

- All models perform bad

657 a woman with short hair indoors

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.053	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



The length of hair is not well modeled

# Case Study

## Hard query

- All models perform bad

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.327	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



Persons **not** in the water in some videos

# Case Study

## Hard query

- All models perform bad

**658 two or more people under a tree**

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.662	0.446	0.795	0.038
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.004	0.001	0.000	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



- Often **only one person or without people**
- Some **in** the tree instead of **under** the tree<sup>26</sup>



# Case Study

## Hard query

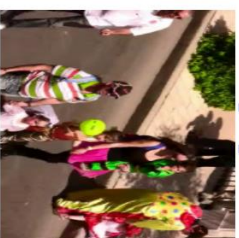
- All models perform bad

**643 people dancing or singing while wearing costumes outdoors**

topic id	run1	run2	run3	run4
641	0.281	0.279	0.232	0.231
642	0.522	0.541	0.454	0.516
643	0.092	0.097	0.053	0.088
644	0.002	0.000	0.000	0.000
645	0.166	0.168	0.153	0.154
646	0.155	0.158	0.132	0.136
647	0.436	0.436	0.400	0.398
648	0.099	0.092	0.105	0.078
649	0.469	0.500	0.546	0.477
650	0.079	0.083	0.057	0.103
651	0.072	0.080	0.001	0.213
652	0.130	0.132	0.123	0.125
653	0.231	0.286	0.045	0.313
654	0.341	0.337	0.320	0.302
655	0.053	0.047	0.013	0.087
656	0.629	0.627	0.618	0.559
657	0.084	0.077	0.068	0.031
658	0.055	0.050	0.080	0.046
659	0.342	0.355	0.323	0.371
660	0.476	0.477	0.512	0.455



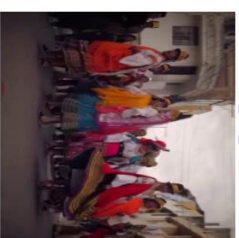
shot07269\_7\_14



shot00164\_34\_23



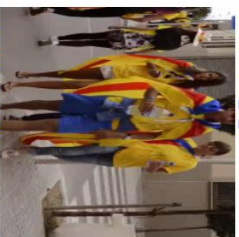
shot00202\_7\_36



shot05029\_366\_0



shot03357\_9\_0



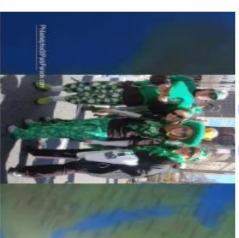
shot04888\_79\_0



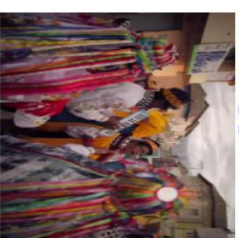
shot00943\_125\_0



shot03967\_79\_87



shot07224\_61\_14



shot05029\_364\_12

**people standing or walking**





# Conclusions

## To boost AVS performance

- Multi-space multi-loss Learning
- Appending extra C3D feature
- Pre-training on image caption dataset
- Late average fusion

## Understanding **fine-grained** queries is still hard

- Attributes: number of persons, length of women's hair, etc.
- Actions: dancing, singing
- Positions: in the water, under a tree



# Reproducibility & Reference

**Code & Resources:** <https://github.com/li-xirong/sea> (in preparation)

**Papers:**

- SEA: Sentence encoder assembly for video retrieval by textual queries. *IEEE Trans. Multimedia 2021*.  
<https://arxiv.org/abs/2011.12091>
- Renmin University of China at TRECVID 2020: Sentence Encoder Assembly for Ad-hoc Video Search, *TRECVID 2020 Workshop*

**Contact:** [xirong@ruc.edu.cn](mailto:xirong@ruc.edu.cn)

# Thanks!